# EDA Assignment - Lending Club Case Study

## Umesh Wagharalkar

## Acknowledgements

- This project was inspired by upGrad

- This project was based on https://www.lendingclub.com/.

## Contact

Created by [UmeshWagharalkar] - feel free to contact me!

# Lending Club Case Study

- In this case study, apart from applying the techniques of EDA, it will also develop a basic understanding of risk analytics in banking and financial services and understand how data is used to minimise the risk of losing money while lending to customers.
- Consumer finance company which specialises in lending various types of loans to urban customers. When the company receives a loan application, the company has to make a decision for loan approval based on the applicant's profile?
- The aim of this case study is to find out the defaulters and non defaulters
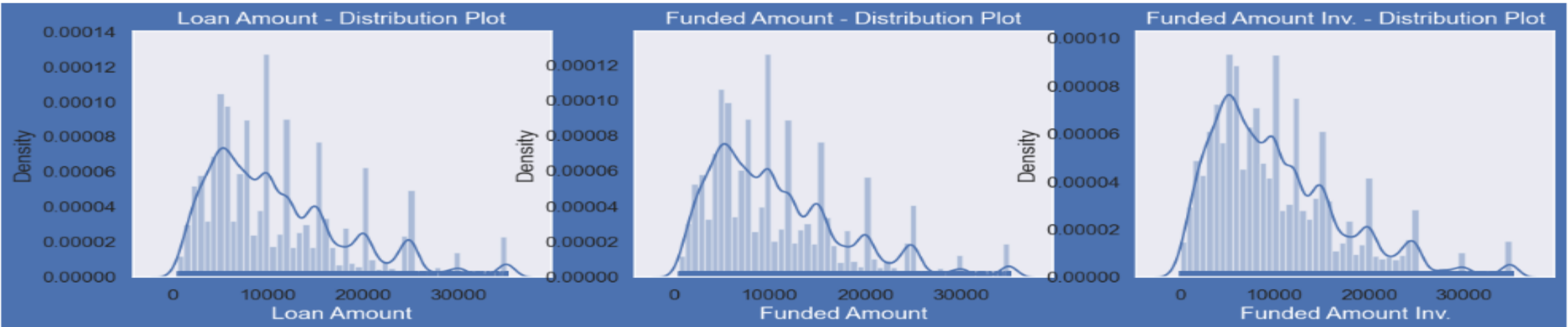- Landing Club data is used for this analysis

## Table of Contents

# Conclusions

- Most of the loans taken for debt consolidation(47%) and Credit card bill payment
- Average intrest rate is 12 %
- Most of the Loan amounts are in range of 5000 - 15000
- Most of the Interest Rates on loans are in range of 10% - 15%
- Most of the borrower's Annual incomes are in range of 40000- 80000
- 14% loans were charged off out of total loan issued
- Most of the loans were taken for the purpose of debt consolidation & paying credit card bill. Number of chraged off count also high too for these loans.
- Most of the applicants are living in rented home or mortgazed their home.
- Loan amount, investor amount, funding amount are strongly correlated.
- Annual income with DTI(Debt-to-income ratio) is negatively correalted.
- Income range 80000+ has less chances of charged off.
- Income range 0-20000 has high chances of charged off
- Small Business applicants have high chnaces of getting charged off.
- Chances of charged off is increasing with grade moving from "A" towards "G"
- Charged off proportion is increasing with higher intrest rates.
- State NE has very high chances of charged off but number of applications are too low.
- States NV,CA and FL states shows good number of charged offs in good number of applications.
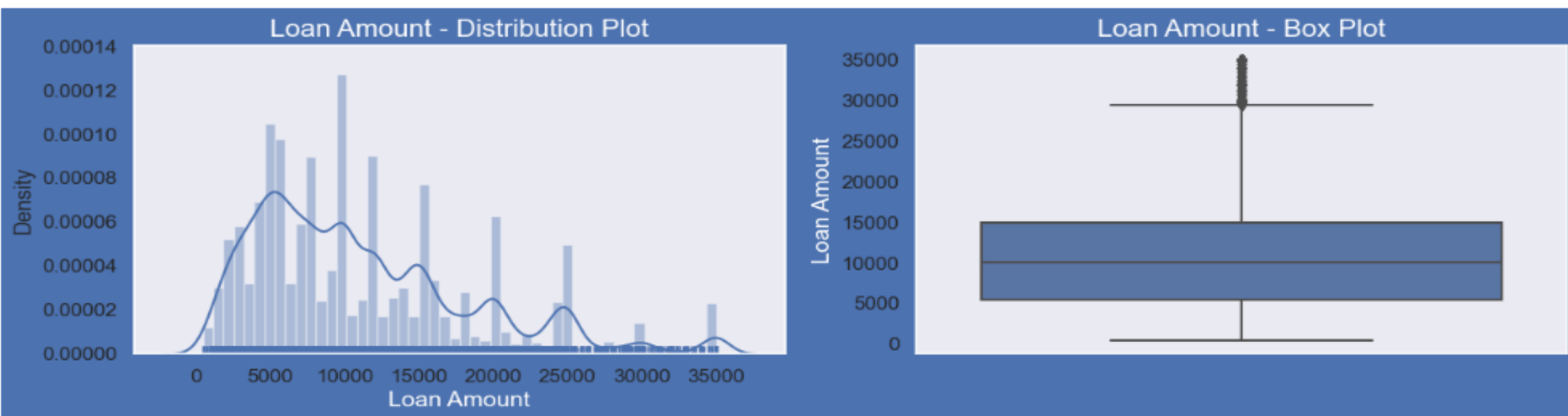
# Univariate Analysis -

```python
# Lets see distribution of three loan amount fields using distribution plot.
# Quantitative Variables

plt.figure(figsize=(15,8),facecolor='b')
sns.set_style("dark")
# subplot 1
plt.subplot(2, 3, 1)
ax = sns.distplot(data['loan_amnt'],rug = True)
ax.set_title('Loan Amount - Distribution Plot',fontsize=14,color='w')
ax.set_xlabel('Loan Amount',fontsize=14,color='w')
# subplot 2
plt.subplot(2, 3, 2)
ax = sns.distplot(data['funded_amnt'],rug = True)
ax.set_title('Funded Amount - Distribution Plot',fontsize=14,color='w')
ax.set_xlabel('Funded Amount',fontsize=14,color='w')
# subplot 2
plt.subplot(2, 3, 3)
ax = sns.distplot(data['funded_amnt_inv'],rug = True)
ax.set_title('Funded Amount Inv. - Distribution Plot',fontsize=14,color='w')
ax.set_xlabel('Funded Amount Inv.',fontsize=14,color='w')
plt.show()

# Observation:
# Distribution of amounts for all three looks very much similar.
# We will work with only loan amount column for rest of our analysis.
```

```python
# Univariate Analysis on Loan amount-Quantitative Variables

plt.figure(figsize=(15,8),facecolor='b')
sns.set_style("dark")
# subplot 1
plt.subplot(2, 2, 1)
ax = sns.distplot(data['loan_amnt'],rug = True)
ax.set_title('Loan Amount - Distribution Plot',fontsize=16,color='w')
ax.set_xlabel('Loan Amount',fontsize=14,color='w')
# subplot 2
plt.subplot(2, 2, 2)
ax = sns.boxplot(y=data['loan_amnt'])
ax.set_title('Loan Amount - Box Plot',fontsize=16,color='w')
ax.set_ylabel('Loan Amount',fontsize=14,color='w')
plt.show()

# Observations :
# Below plots show that most of the Loan amounts are in range of 5000 - 15000
```
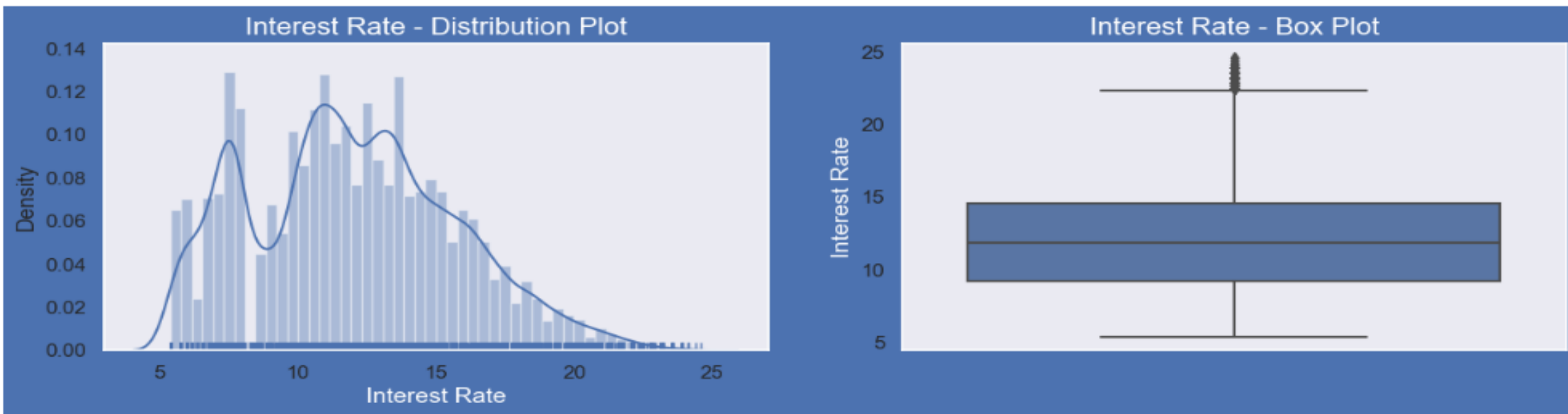
```
1   # Univariate Analysis on Intrest Rate-Quantitative Variables
2
3   plt.figure(figsize=(15,8),facecolor='b')
4   sns.set_style("dark")
5   # subplot 1
6   plt.subplot(2, 2, 1)
7   ax = sns.distplot(data['int_rate'],rug = True)
8   ax.set_title('Interest Rate - Distribution Plot',fontsize=16,color='w')
9   ax.set_xlabel('Interest Rate',fontsize=14,color='w')
10  # subplot 2
11  plt.subplot(2, 2, 2)
12  ax = sns.boxplot(y=data['int_rate'])
13  ax.set_title('Interest Rate - Box Plot',fontsize=16,color='w')
14  ax.set_ylabel('Interest Rate',fontsize=14,color='w')
15  plt.show()
16
17  # Observations :
18  # Below plots show that most of the Interest Rates on Loans are in range of 10% - 15%
```
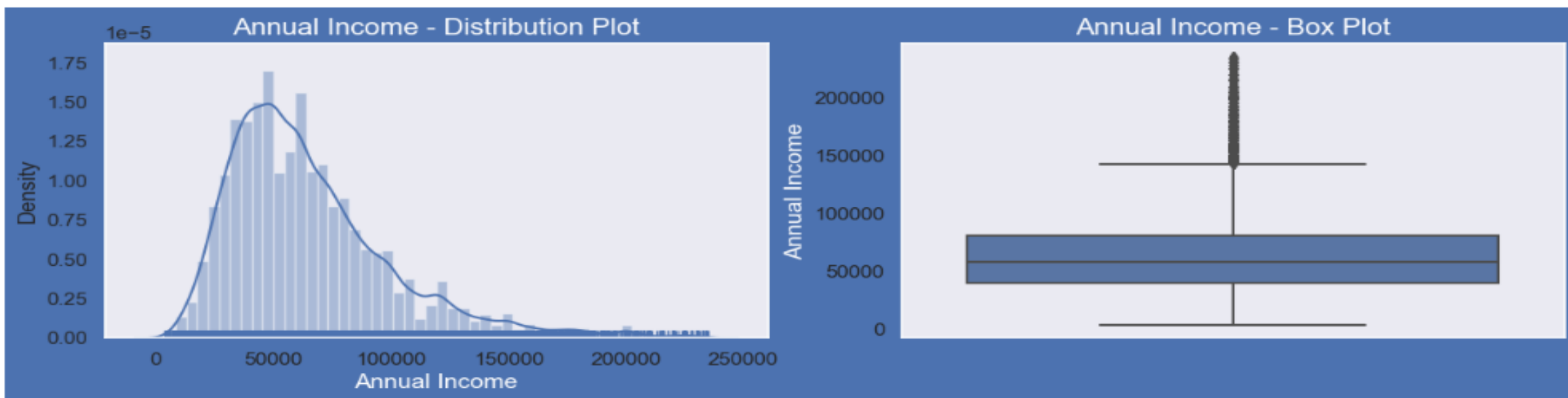
```python
# Univariate Analysis on Annual Income - Quantitative Variables

plt.figure(figsize=(15,8),facecolor='b')
sns.set_style("dark")
# subplot 1
plt.subplot(2, 2, 1)
ax = sns.distplot(data['annual_inc'],rug = True)
ax.set_title('Annual Income - Distribution Plot',fontsize=16,color='w')
ax.set_xlabel('Annual Income',fontsize=14,color='w')
# subplot 2
plt.subplot(2, 2, 2)
plt.title('Annual Income - Box Plot')
ax = sns.boxplot(y=data['annual_inc'])
ax.set_title('Annual Income - Box Plot',fontsize=16,color='w')
ax.set_ylabel('Annual Income',fontsize=14,color='w')
plt.show()

# Observations :
# Below plots show that most of the borrower's Annual incomes are in range of 40000- 80000
```
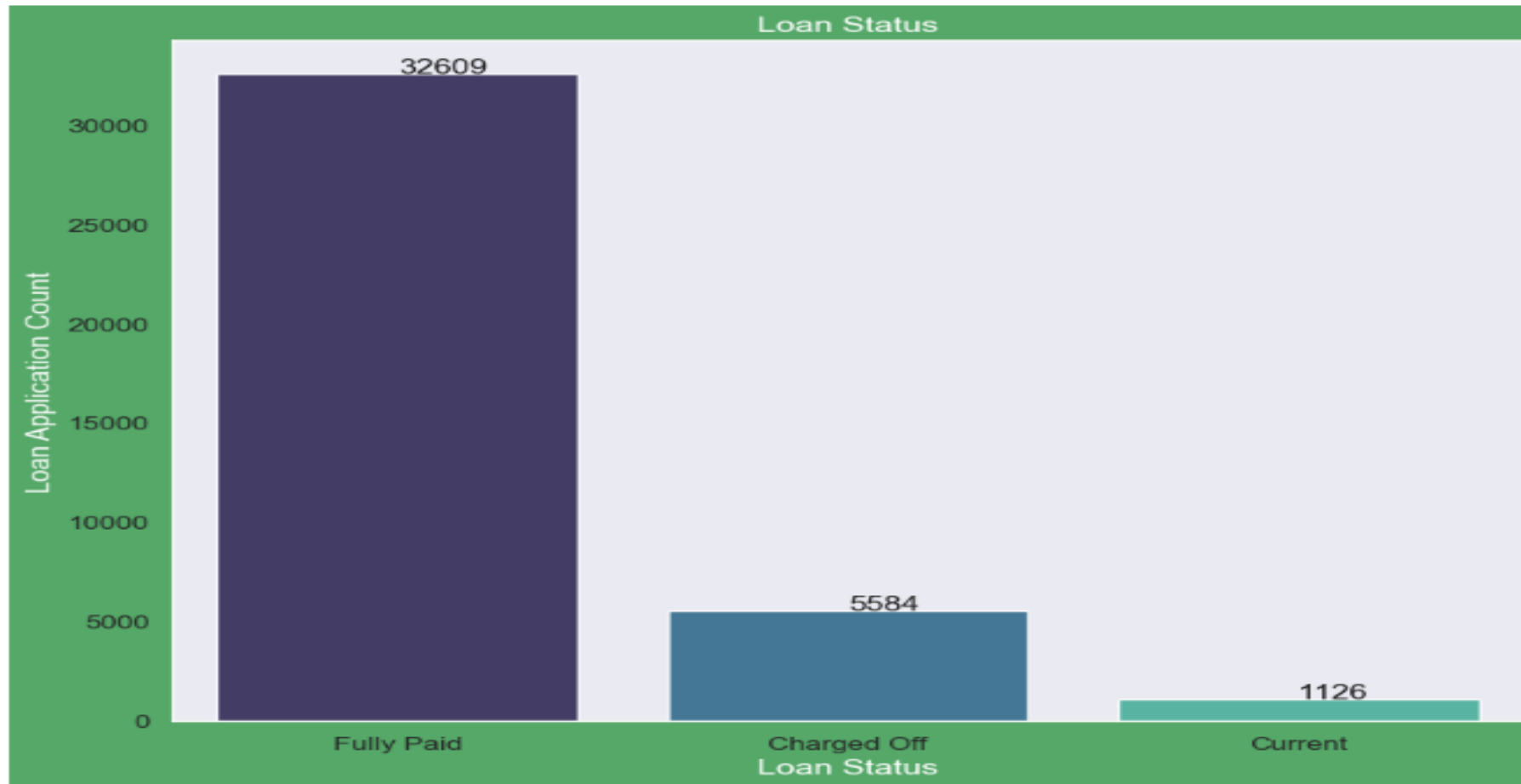
```
1  # Univariate Analysis - Unordered Categorical Variables - Loan Status
2
3  plt.figure(figsize=(10,8),facecolor='g')
4  sns.set_style("dark")
5  ax = sns.countplot(x="loan_status",data=data,palette='mako')
6  ax.set_title('Loan Status',fontsize=14,color='w')
7  ax.set_xlabel('Loan Status',fontsize=14,color = 'w')
8  ax.set_ylabel('Loan Application Count',fontsize=14,color = 'w')
9  # To show count of values above bars
10 s=data['loan_status'].value_counts()
11 for i, v in s.reset_index().iterrows():
12     ax.text(i, v.loan_status + 0.3 , v.loan_status, color='k')
13
14 # Observations :
15 # Below plot shows that close to 14% Loans were charged off out of total Loan issued.
```
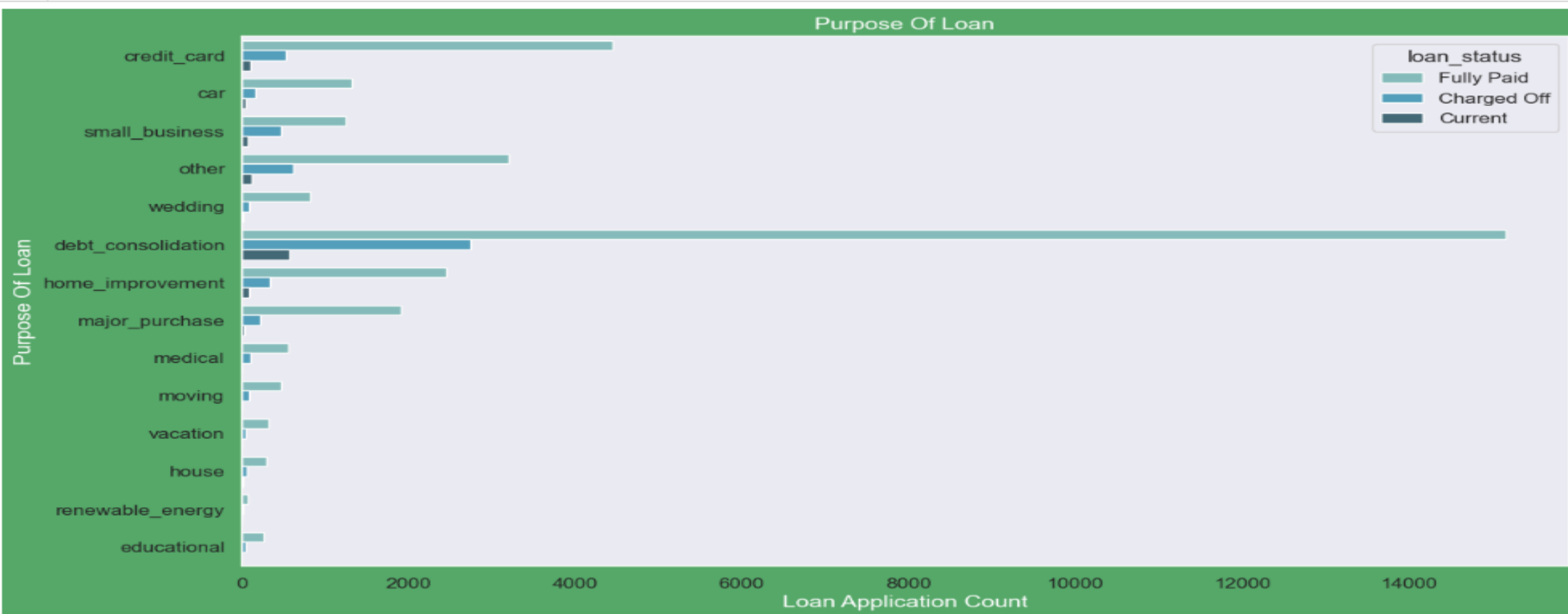
```
1   # Univariate Analysis - Unordered Categorical Variables - Purpose Of Loan
2
3   plt.figure(figsize=(14,8),facecolor='g')
4   sns.set_style("dark")
5   ax = sns.countplot(y="purpose",data=data,hue='loan_status',palette='GnBu_d')
6   ax.set_title('Purpose Of Loan',fontsize=14,color='w')
7   ax.set_ylabel('Purpose Of Loan',fontsize=14,color = 'w')
8   ax.set_xlabel('Loan Application Count',fontsize=14,color = 'w')
9   plt.show()
10
11  # Observations :
12  # Below plot shows that most of the loans were taken for the purpose of debt consolidation & paying credit card bill.
13  # Number of chraged off count also high too for these loans.
```
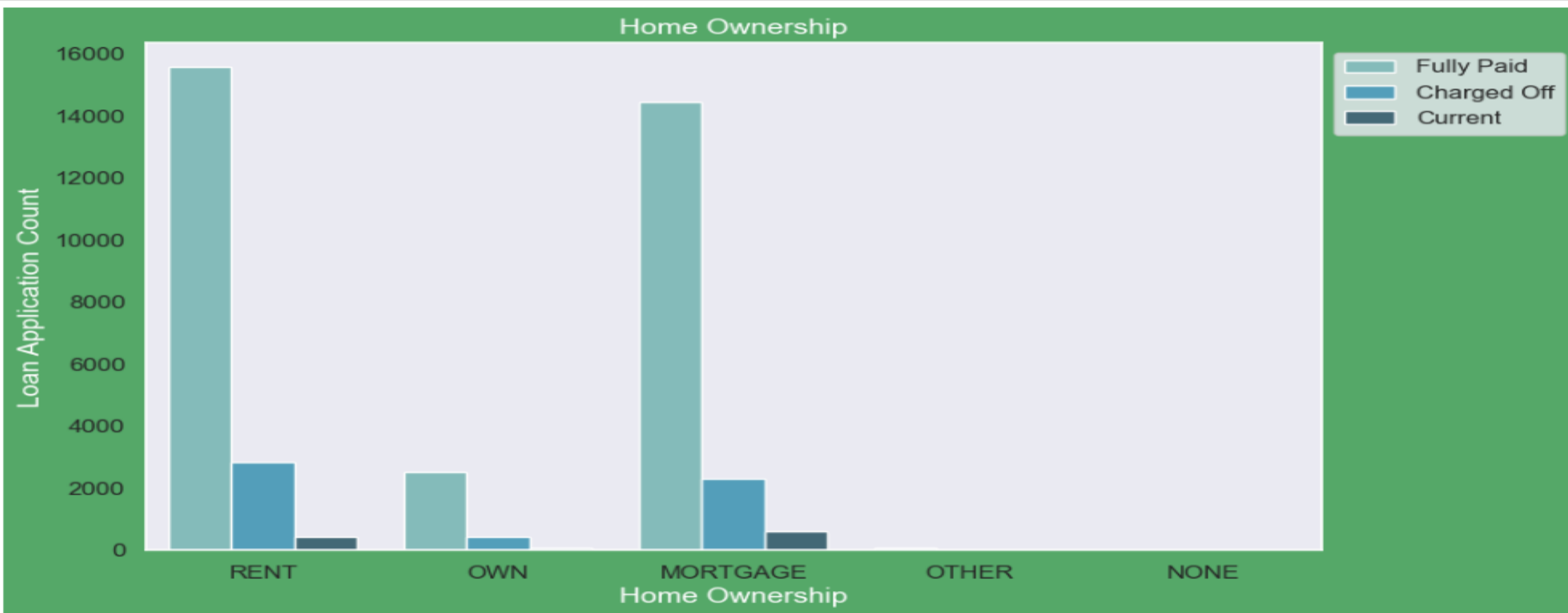


Purpose Of Loan

```
# Univariate Analysis - Unordered Categorical Variables - Home Ownership

plt.figure(figsize=(10,6),facecolor='g')
ax = sns.countplot(x="home_ownership",data=data,hue='loan_status',palette='GnBu_d')
ax.legend(bbox_to_anchor=(1, 1))
ax.set_title('Home Ownership',fontsize=14,color='w')
ax.set_xlabel('Home Ownership',fontsize=14,color = 'w')
ax.set_ylabel('Loan Application Count',fontsize=14,color = 'w')
plt.show()

# Observations :
# Below plot shows that most of them living in rented home or mortgazed their home.
# Applicant numbers are high from these categories so charged off is high too.
```
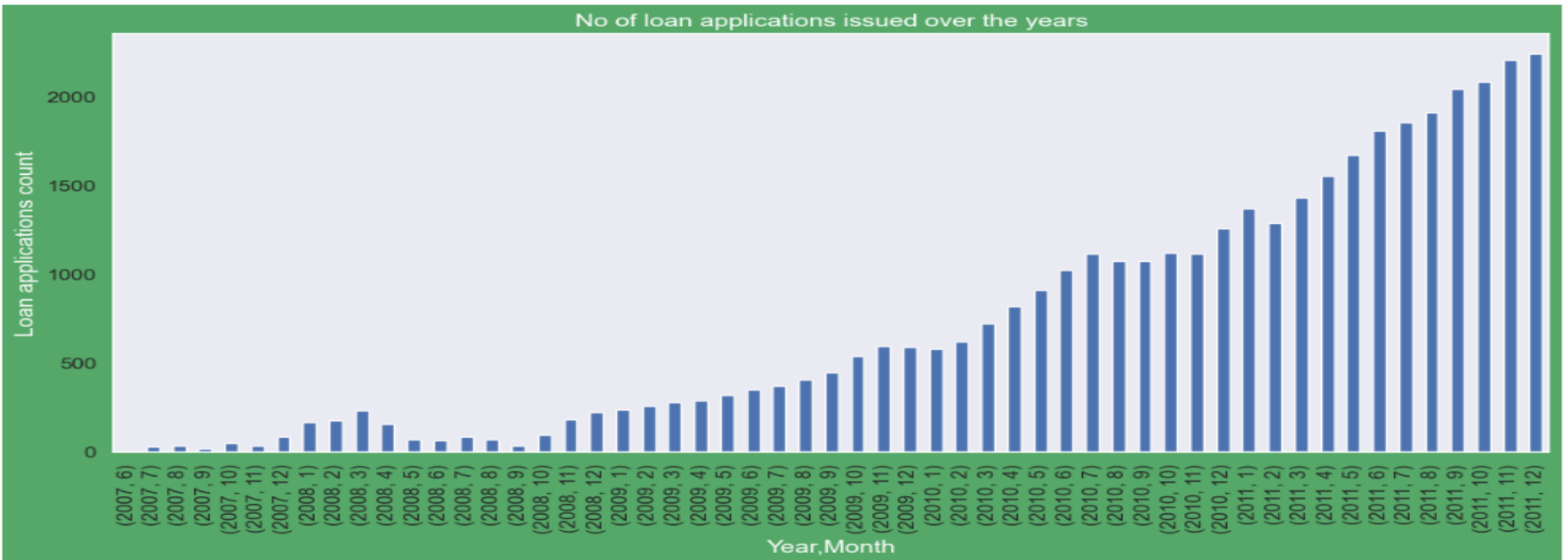
```
1   # Derived Column - Ordered Categorical Variables
2   # Let us look into number of Loans which were approved every year/month
3   # Lets use derived column year to check pattern of loan issuing over the years.
4   plt.figure(figsize=(14,6),facecolor='g')
5   data.groupby(['year','month']).id.count().plot(kind='bar')
6   plt.ylabel('Loan applications count',fontsize=14,color='w')
7   plt.xlabel('Year,Month',fontsize=14,color = 'w')
8   plt.title("No of loan applications issued over the years",fontsize=14,color='w')
9   plt.show()
10
11
12  # Observation is that count of Loan application is increasing every passing year.
13  # so increase in number of Loan applications are adding more to number of charged off applications.
14  # number of Loans issued in 2008( May-October) got dipped, may be due to Recession.
```
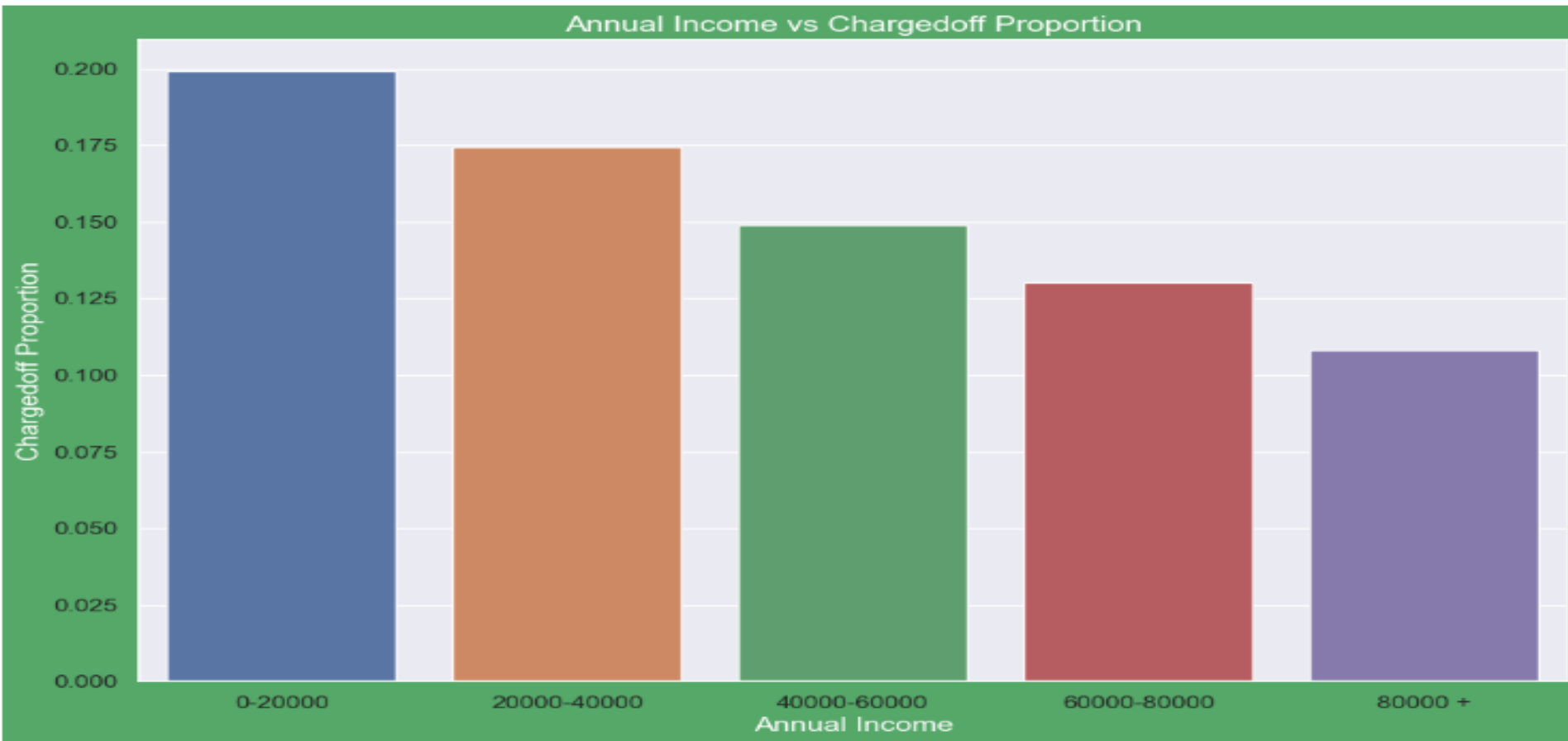


No of loan applications issued over the years

```python
# Drawing bar plots on data calculated above. Try to visualize the pattern to understand the data better.

fig, ax1 = plt.subplots(figsize=(12, 8),facecolor='g')
ax1.set_title('Annual Income vs Chargedoff Proportion',fontsize=15,color = 'w')
ax1=sns.barplot(x='annual_inc_cats', y='Chargedoff_Proportion', data=inc_range_vs_loan)
ax1.set_ylabel('Chargedoff Proportion',fontsize=14,color = 'w')
ax1.set_xlabel('Annual Income',fontsize=14,color='w')
plt.show()

# Observations:
# Income range 80000+  has Less chances of charged off.
# Income range 0-20000 has high chances of charged off.
# Notice that with increase in annual income charged off proportion got decreased.
```
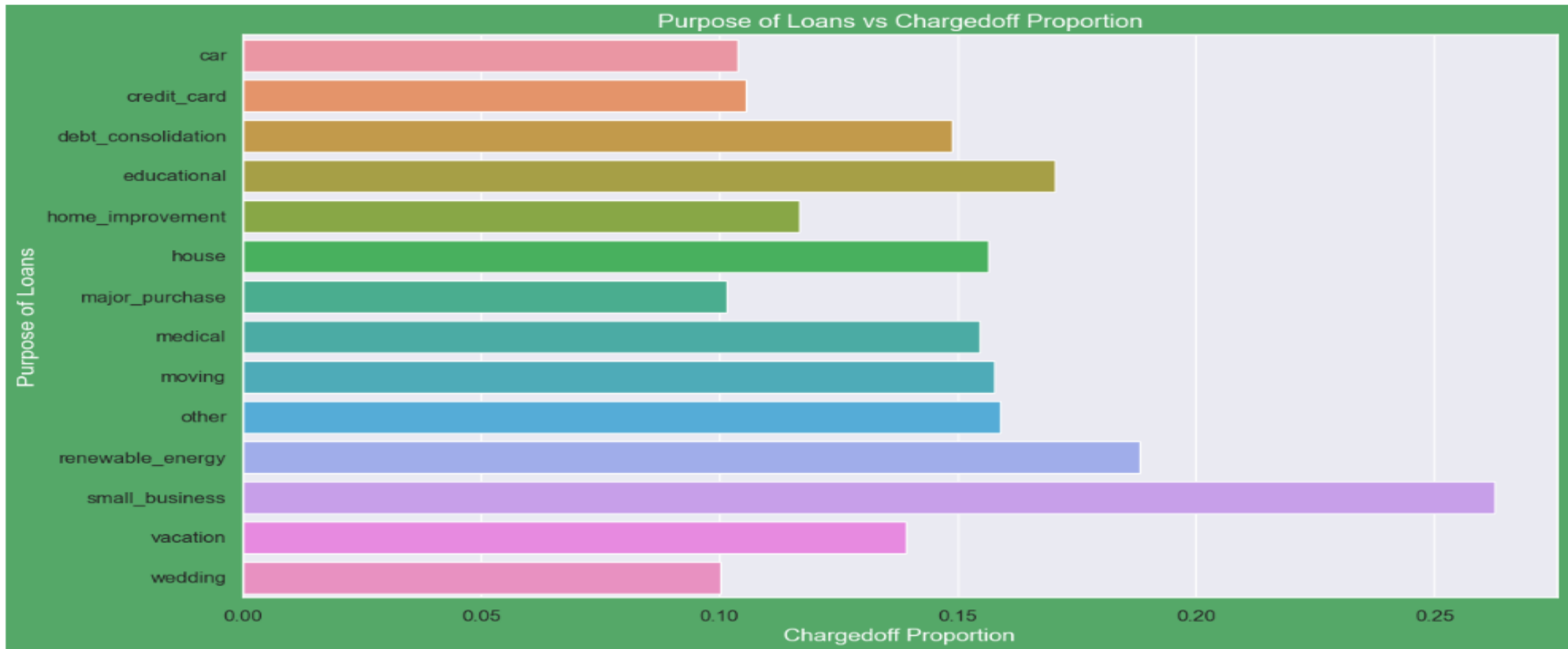


Annual Income vs Chargedoff Proportion

```
1   # Drowing bar plots on data calculated above. Try to visualize the pattern to understand the data better.
2   # Pairs of continuous variables.
3   fig, ax1 = plt.subplots(figsize=(14, 8),facecolor='g')
4   ax1.set_title('Purpose of Loans vs Chargedoff Proportion',fontsize=15,color = 'w')
5   ax1=sns.barplot(y='purpose', x='Chargedoff_Proportion', data=purpose_vs_loan)
6   ax1.set_ylabel('Purpose of Loans',fontsize=14,color='w')
7   ax1.set_xlabel('Chargedoff Proportion',fontsize=14,color = 'w')
8   plt.show()
9
10  # Observations:
11  # small Business applicants have high chances of getting charged off.
12  # renewable_energy where chanrged off proportion is better as compare to other categories.
```
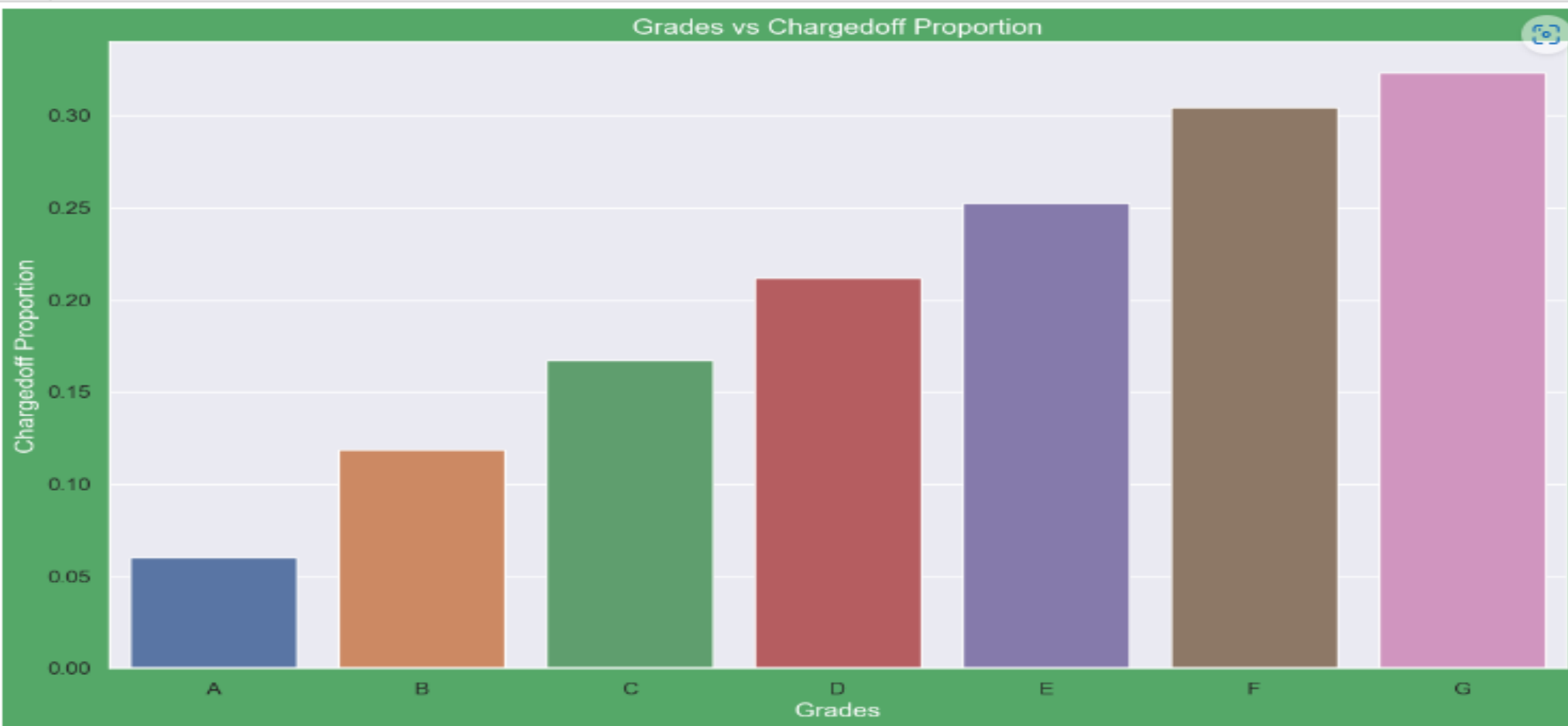
```
1  # Drawing bar plots on data calculated above. Try to visualize the pattern to understand the data better.
2
3  fig, ax1 = plt.subplots(figsize=(14, 8),facecolor='g')
4  ax1.set_title('Grades vs Chargedoff Proportion',fontsize=15,color='w')
5  ax1=sns.barplot(x='grade', y='Chargedoff_Proportion', data=grade_vs_loan)
6  ax1.set_xlabel('Grades',fontsize=14,color='w')
7  ax1.set_ylabel('Chargedoff Proportion',fontsize=14,color ='w')
8  plt.show()
9
10 # Observations:
11 # Grade "A" has very Less chances of charged off.
12 # Grade "F" and "G" have very high chances of charged off.
13 # Chances of charged of is increasing with grade moving from "A" towards "G"
```
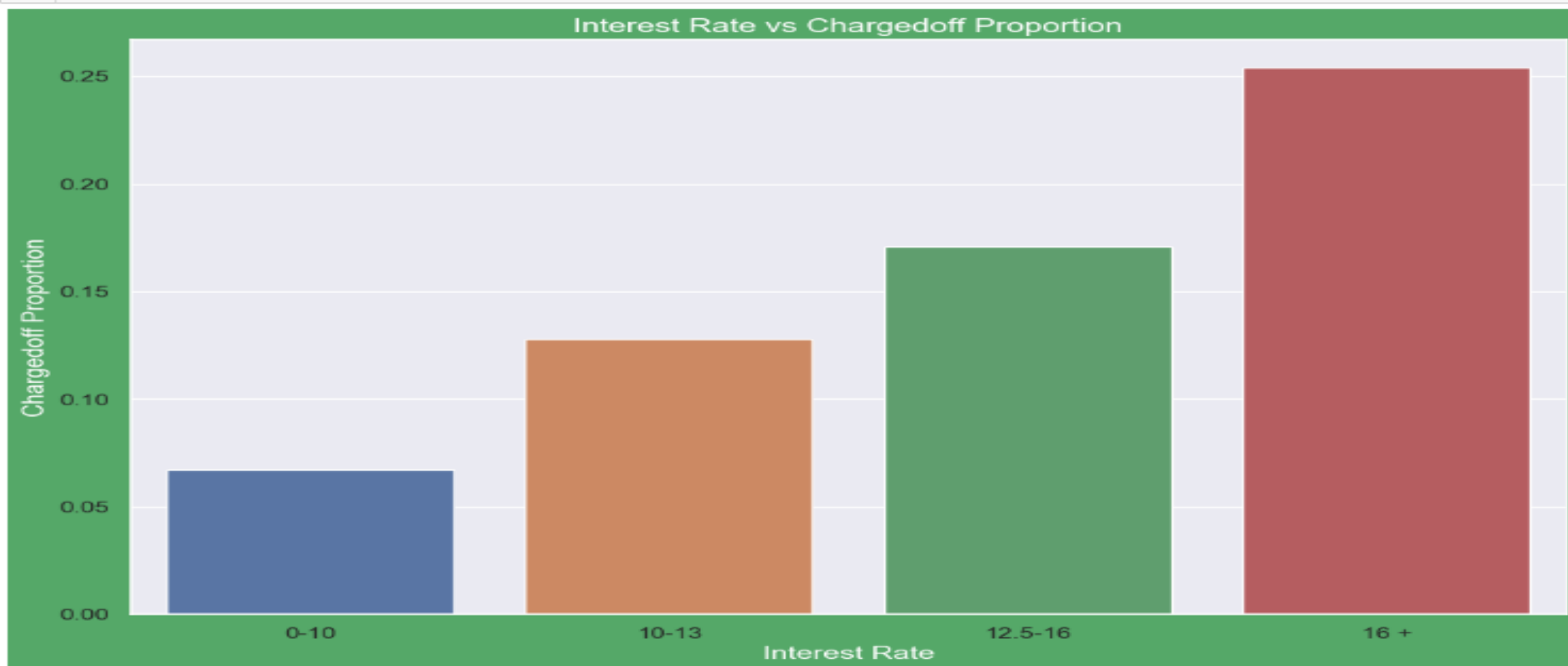


Grades vs Chargedoff Proportion

```
1  # Drawing some bar plots on data calculated above. Try to visualize the pattern to understand the data better.
2
3  fig, ax1 = plt.subplots(figsize=(12, 8),facecolor='g')
4  ax1.set_title('Interest Rate vs Chargedoff Proportion',fontsize=15,color='w')
5  ax1=sns.barplot(x='int_rate_cats', y='Chargedoff_Proportion', data=interest_vs_loan)
6  ax1.set_xlabel('Interest Rate',fontsize=14,color='w')
7  ax1.set_ylabel('Chargedoff Proportion',fontsize=14,color = 'w')
8  plt.show()
9
10 # Observations:
11 # interest rate Less than 10% has very Less chances of charged off. Intrest rates are starting from minimin 5 %.
12 # interest rate more than 16% has good chnaces of charged off as compared to other category intrest rates.
13 # Charged off proportion is increasing with higher intrest rates.
```
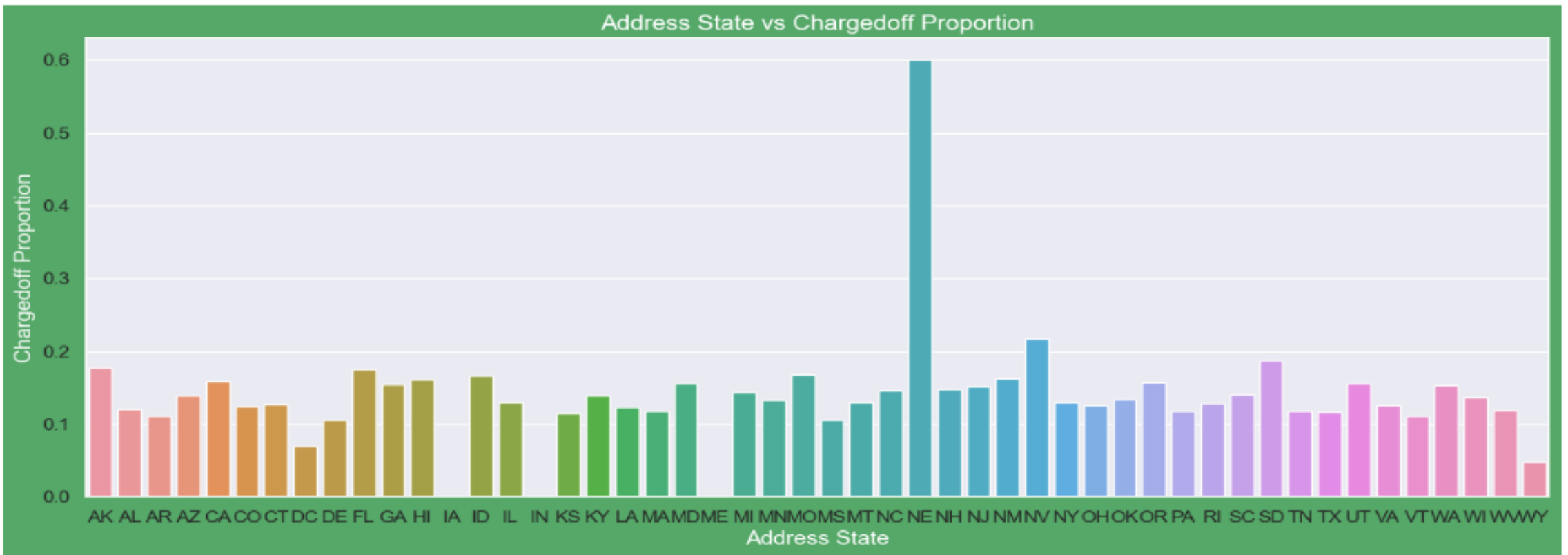


Interest Rate vs Chargedoff Proportion

```
 1  # Drawing bar plots on data calculated above. Try to visualize the pattern to understand the data better.
 2
 3  fig, ax1 = plt.subplots(figsize=(16, 6),facecolor='g')
 4  ax1.set_title('Address State vs Chargedoff Proportion',fontsize=15,color='w')
 5  ax1=sns.barplot(x='addr_state', y='Chargedoff_Proportion', data=state_vs_loan)
 6  ax1.set_xlabel('Address State',fontsize=14,color='w')
 7  ax1.set_ylabel('Chargedoff Proportion',fontsize=14,color = 'w')
 8  plt.show()
 9
10  # Observations:
11  # states NE has very high chances of charged off but number of applications are too low to make any decisions.
12  # NV,CA and FL states shows good number of charged offs in good number of applications.
```
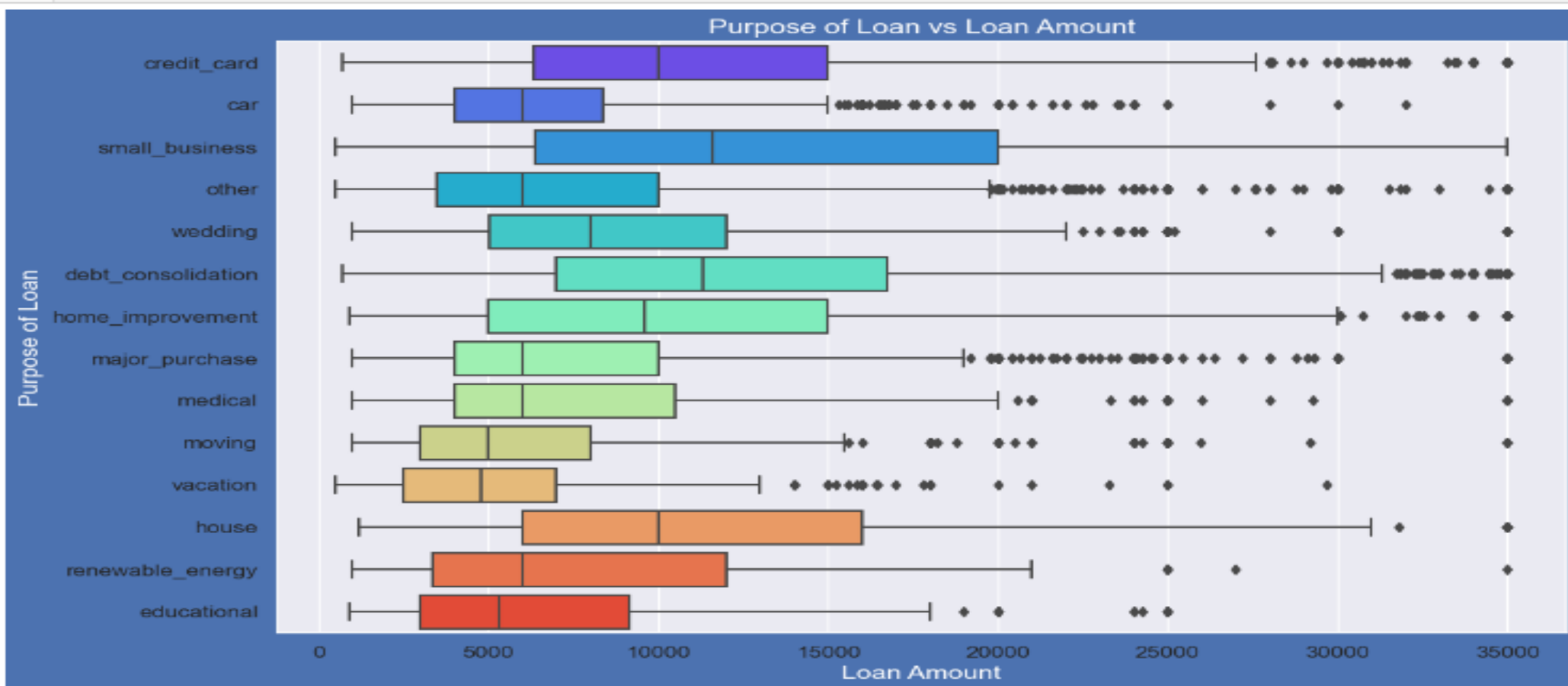


Address State vs Chargedoff Proportion

```
1  # Bivariate Analysis - Puprose of Loan vs Loan amount
2  # Box Plot
3
4  plt.figure(figsize=(12,8),facecolor='b')
5  ax = sns.boxplot(y='purpose', x='loan_amnt', data =data,palette='rainbow')
6  ax.set_title('Purpose of Loan vs Loan Amount',fontsize=15,color='w')
7  ax.set_ylabel('Purpose of Loan',fontsize=14,color = 'w')
8  ax.set_xlabel('Loan Amount',fontsize=14,color = 'w')
9  plt.show()
10
11 # Observations:
12 # Median,95th percentile,75th percentile of Loan amount is highest for Loan taken for small business purpose among all purpo
13 # Debt consolidation is second and Credit card comes 3rd.
```
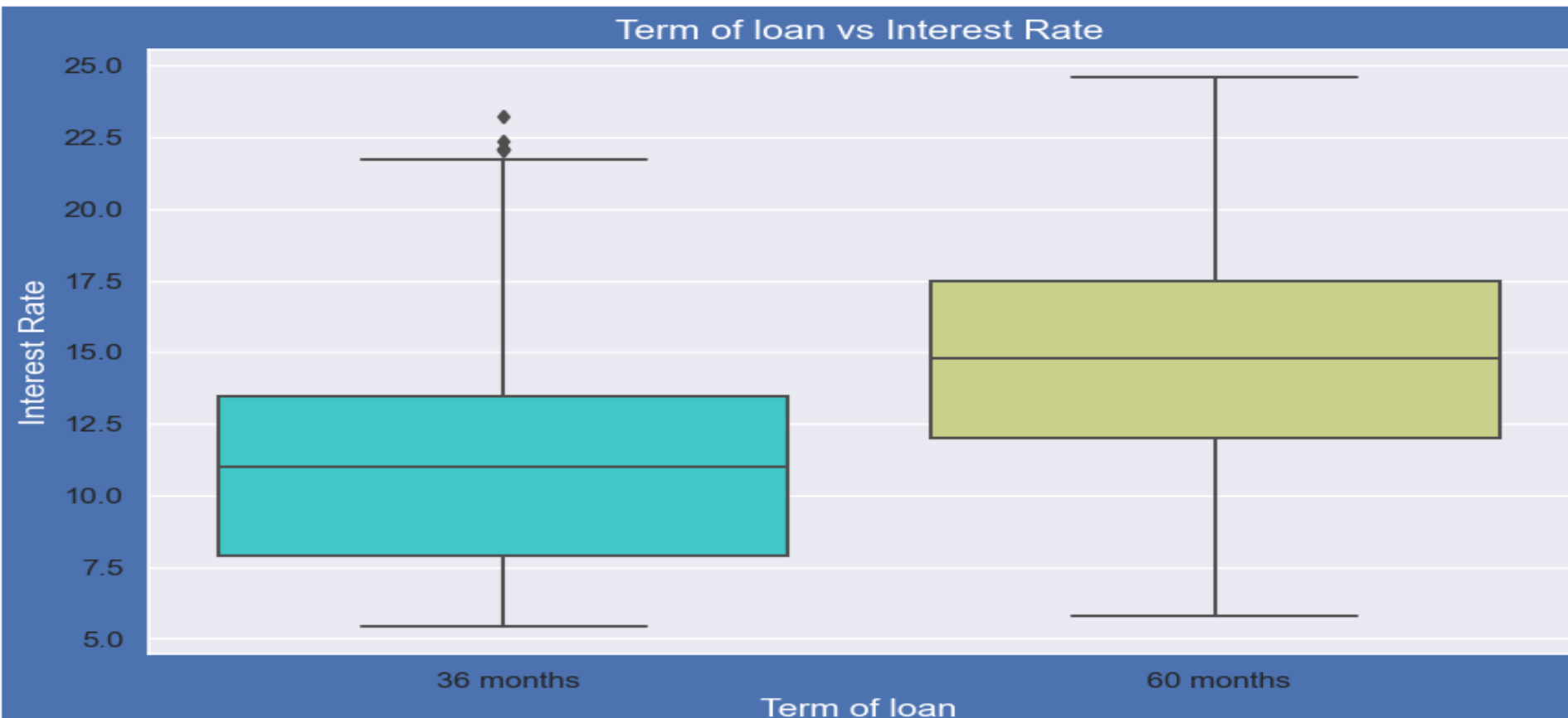


Purpose of Loan vs Loan Amount

```
1   # Bivariate Analysis - Term of loan vs Interest Rate
2   # Box Plot
3
4   plt.figure(figsize=(10,6),facecolor='b')
5   ax = sns.boxplot(y='int_rate', x='term', data =data,palette='rainbow')
6   ax.set_title('Term of loan vs Interest Rate',fontsize=15,color='w')
7   ax.set_ylabel('Interest Rate',fontsize=14,color = 'w')
8   ax.set_xlabel('Term of loan',fontsize=14,color = 'w')
9   plt.show()
10
11  # Observations:
12  # It is clear that avearge intrest rate is higher for 60 months Loan term.
13  # Most of the Loans issued for Longer term had higher intrest rates for repayement.
```
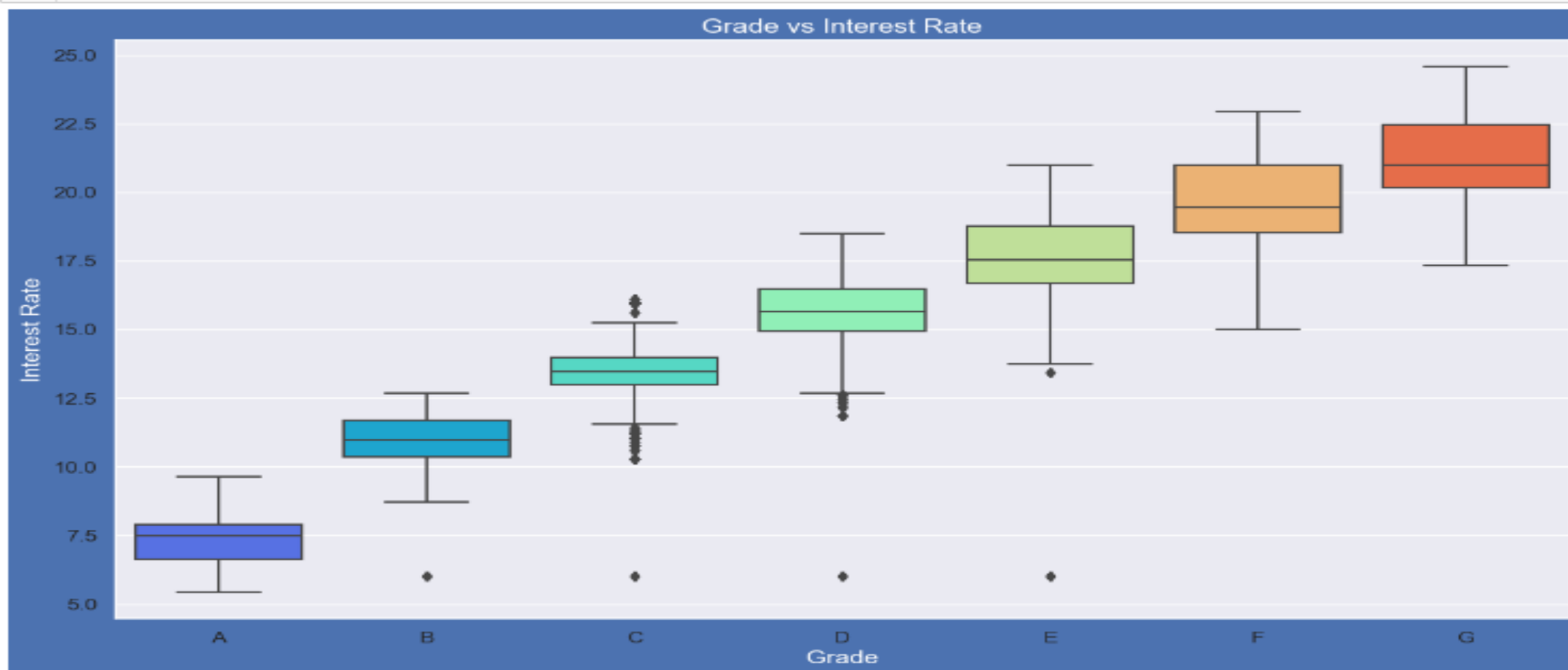


Term of loan vs Interest Rate

```python
# Bivariate Analysis - Grade vs Interest Rate
# Box Plot

plt.figure(figsize=(14,8),facecolor='b')
ax = sns.boxplot(y='int_rate', x='grade', data =data,palette='rainbow',order = 'ABCDEFG')
ax.set_title('Grade vs Interest Rate',fontsize=15,color='w')
ax.set_ylabel('Interest Rate',fontsize=14,color = 'w')
ax.set_xlabel('Grade',fontsize=14,color = 'w')
plt.show()

# Observations:
# A-grade is a top Letter grade for a Lender to assign to a borrower.
# The higher the borrower's credit grade,the Lower the interest rate offered to that borrower on a Loan.
# It is clear that intrest rate is increasing with grades moving from A to F.
```
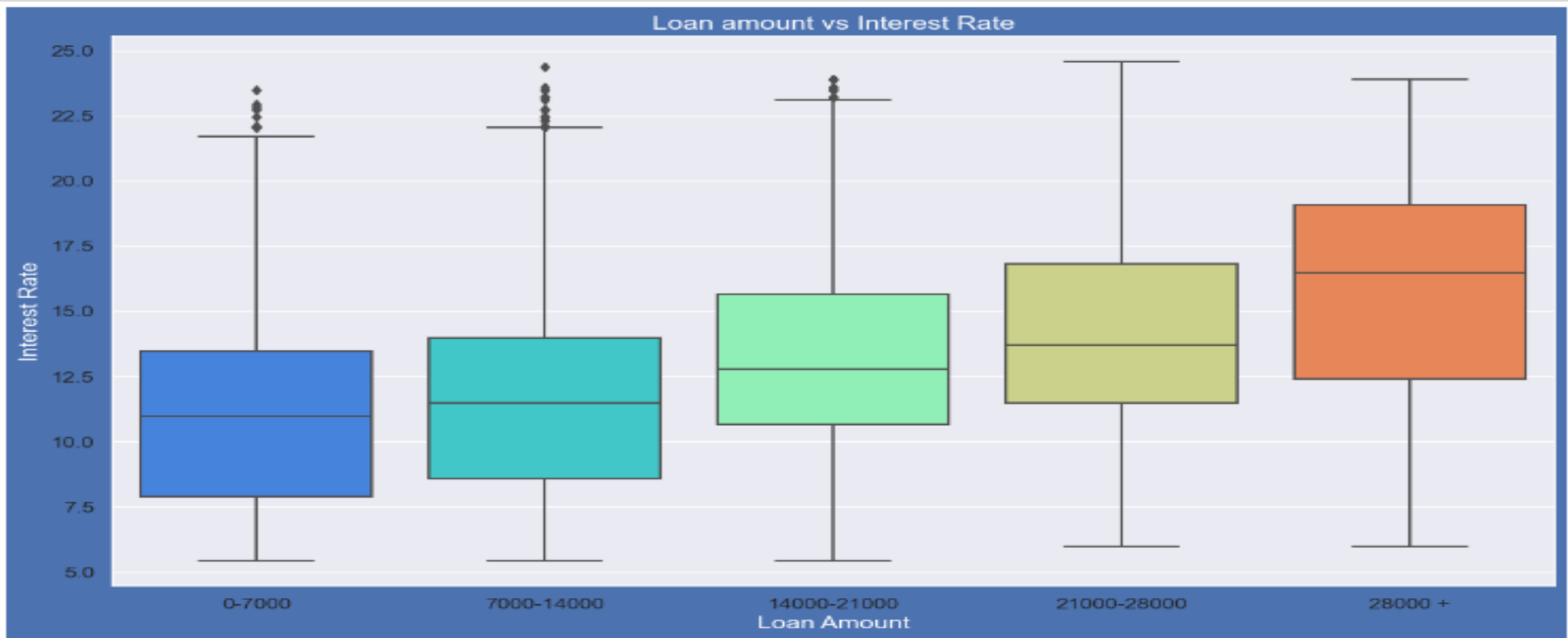


Grade vs Interest Rate

```
1   # Bivariate Analysis - Loan Amount vs Interest Rate
2   # Box Plot
3
4   plt.figure(figsize=(14,8),facecolor='b')
5   ax = sns.boxplot(y='int_rate', x='loan_amnt_cats', data =data,palette='rainbow')
6   ax.set_title('Loan amount vs Interest Rate',fontsize=15,color='w')
7   ax.set_ylabel('Interest Rate',fontsize=14,color = 'w')
8   ax.set_xlabel('Loan Amount',fontsize=14,color = 'w')
9   plt.show()
10
11  # Observations:
12  # It is clear that intrest rate is increasing with Loan amount increase.
13  # probably when Loan amount is more it is taken for Longer Loan term, we saw earlier that Longer the Loan term more the
14  # interest rate.
```
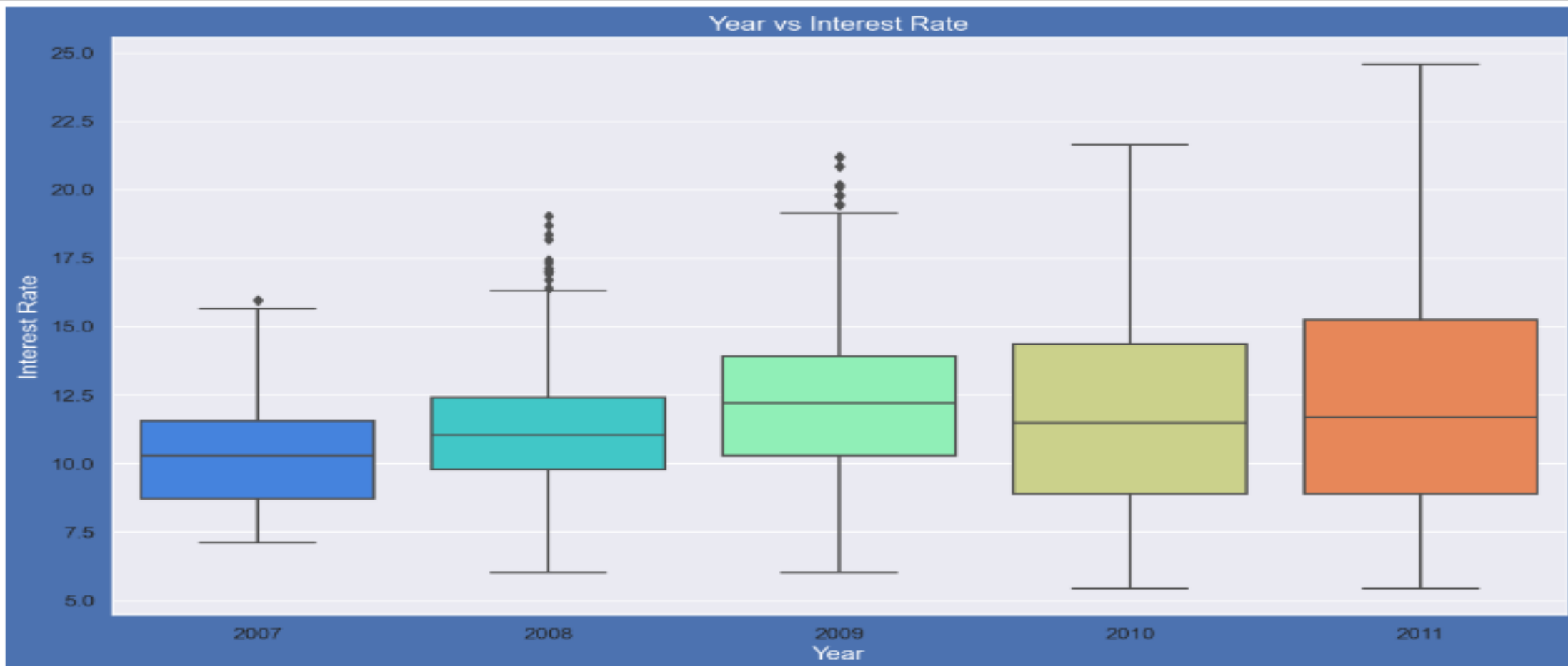


Loan amount vs Interest Rate

```
1  # Bivariate Analysis - year vs Interest Rate
2  # Box Plot
3
4  plt.figure(figsize=(14,8),facecolor='b')
5  ax = sns.boxplot(y='int_rate', x='year', data =data,palette='rainbow')
6  ax.set_title('Year vs Interest Rate',fontsize=15,color='w')
7  ax.set_ylabel('Interest Rate',fontsize=14,color = 'w')
8  ax.set_xlabel('Year',fontsize=14,color = 'w')
9  plt.show()
10
11 # Observations:
12 # Plot shows intrest rate is increasing slowly with increase in year.
```



Year vs Interest Rate

# End of Project